

Purdue University

Purdue e-Pubs

Department of Computer Science Technical
Reports

Department of Computer Science

1973

Modeling A Large Online File System

John Pomeranz

Report Number:

73-108

Pomeranz, John, "Modeling A Large Online File System" (1973). *Department of Computer Science Technical Reports*. Paper 60.
<https://docs.lib.purdue.edu/cstech/60>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries.
Please contact epubs@purdue.edu for additional information.

MODELING A LARGE ONLINE FILE SYSTEM

John Pomeranz

**Computer Science Department
Purdue University**

**March 1973
CSD TR 108**

MODELING A LARGE ONLINE FILE SYSTEM

INTRODUCTION

The data management system in use at Purdue University's computing center supports both batch and interactive users. To determine a strategy for disk storage allocation, the individual file sizes of the two user groups are separated using hyperexponential curve fitting with linear regression techniques. The techniques are applicable when the exponential parameters of the hyperexponential distribution differ by an order of magnitude.

Composition of Exponential Distributions

The exponential and hyperexponential distributions are the two simplest examples of composition of exponential distributions. Any distribution in this family is a weighted average of exponential distributions. Let $\lambda_1, \dots, \lambda_k$ and $\omega_1, \dots, \omega_k$ be positive constants with $\omega_1 + \dots + \omega_k = 1$. A random variable with probability density

$$(1) \quad f(x) = \sum_{i=1}^k \omega_i \lambda_i e^{-\lambda_i x}$$

and cumulative distribution

$$\begin{aligned} (2) \quad F(x) &= \sum_{i=1}^k \omega_i (1 - e^{-\lambda_i x}) \\ &= 1 - \sum_{i=1}^k \omega_i e^{-\lambda_i x} \end{aligned}$$

is a k-fold exponential distribution.

It is generally quite difficult to identify the weights ω_i and parameters λ_i in actual experimental situations, but for $k = 1, 2$ the techniques below are useful and simple.

Linear Regression Model

Let $(x_1, y_1), \dots, (x_n, y_n)$, $n > 1$, be data points for which a straight line fit is sought, i.e., find constants m and b so that the line

$$(3) \quad y = mx + b$$

lies "close to" the n data points. Traditionally, "close to" is defined as the line that can

$$(4) \quad \text{minimize} \quad \sum_{i=1}^n (y_i - (mx_i + b))^2.$$

The well-known formulas satisfying (4) are given [in Freund et al., p. 40] as

$$(5) \quad m = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n y_i \right) \left(\sum_{i=1}^n x_i \right)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2}$$

and

$$(6) \quad b = \frac{1}{n} \left(\sum_{i=1}^n y_i - m \sum_{i=1}^n x_i \right).$$

The predictive value of the model is in equation (3) which characterizes the process that generates y values from x values.

Linear Regression and the Exponential Density

Suppose the n data points are assumed to be from an exponential density with unknown parameter λ . Thus, the curve being sought is of the form

$$(7) \quad y = f(x) = \lambda e^{-\lambda x}.$$

To employ linear regression, the natural logarithms of the y_i are used. The linear regression model is applied to the pairs $(x_i, \ln y_i)$, $i = 1, \dots, n$.

The m and b obtained satisfy

$$\ln y = mx + b$$

or

$$\begin{aligned} (8) \quad y &= e^{mx + b} \\ &= e^b e^{mx} \end{aligned}$$

Note that the criterion satisfied in (4) is now interpreted as a weighted product of ratios between the observed and predicted points. Furthermore, if y is exponentially distributed then it follows that $-m = e^b = \lambda$. The ratio m/e^b is a measure of the validity of the assumption about the distribution.

Alternate Linear Regression for an Exponential Distribution

The method above can be applied to the cumulative exponential distribution as well as the exponential density. Suppose the pairs (x_i, y_i) are ordered so that $x_{i-1} < x_i$, $i = 2, \dots, n$. Let $Y_i = \sum_{j=1}^i y_j$ be the sample distribution function and let $y_i^* = \ln(1 - Y_i)$. Apply the linear regression model to the pairs (x_i, y_i^*) . Using $Y(x)$ to represent the distribution curve, the m and b obtained from the model satisfy

$$y^* = \ln(1 - Y(x)) = mx + b$$

or

$$\begin{aligned} (9) \quad Y(x) &= 1 - e^{mx + b} \\ &= 1 - e^b e^{mx} \end{aligned}$$

When Y is an exponential distribution then $b = 0$ and $-m = \lambda$. The proximity of b to zero is a measure of the validity of the assumption about the distribution.

Linear Regression and the Hyperexponential Distribution

Here the curve sought has density

$$f(x) = \omega \lambda_1 e^{-\lambda_1 x} + (1 - \omega) \lambda_2 e^{-\lambda_2 x}$$

and distribution

$$(10) \quad F(x) = 1 - (\omega e^{-\lambda_1 x} + (1 - \omega) e^{-\lambda_2 x})$$

Transposing in (10) yields

$$(11) \quad 1 - F(x) = \omega e^{-\lambda_1 x} + (1 - \omega) e^{-\lambda_2 x}$$

For very small or very large values of x one exponential component predominates. Suppose $\lambda_1 < \lambda_2$, and for large values of x , $e^{-\lambda_2 x} \approx 0$. Thus (11) can be approximated by

$$(12) \quad 1 - F(x) = \omega e^{-\lambda_1 x}, \quad x \gg 0,$$

and the analysis for the exponential (above) can be applied to obtain $\hat{\omega}$ and $\hat{\lambda}_1$ as estimates solving (12). Substituting $\hat{\omega}$ and $\hat{\lambda}_1$ in (11) yields

$$(13) \quad 1 - F(x) - \hat{\omega} e^{-\hat{\lambda}_1 x} = (1 - \omega) e^{-\lambda_2 x}$$

as a distribution from which $(1 - \omega)$ and λ_2 can be estimated.

Note that when the regression model is applied to (13), a second estimate for ω is obtained that can be used as a check on the overall distribution.

Verifying the Hyperexponential Parameters

If λ_1 and λ_2 are known (or estimated), ω can be found (or verified for consistency) by rewriting (11) as

$$(14) \frac{1 - F(x) - e^{-\lambda_2 x}}{e^{-\lambda_1 x} - e^{-\lambda_2 x}} = \omega$$

and applying the linear regression model to the pairs

$$(x_1, (1 - F(x) - \exp(-\lambda_2 x)) / (\exp(-\lambda_1 x) - \exp(-\lambda_2 x))).$$

The resulting m should be zero and $\omega = b$. Note that regardless of the complexity of (14), ω is the constant term of the linear regression equation.

The Weibull Distribution

A random variable with probability density

$$(15) f(x) = \lambda c x^{c-1} e^{-\lambda x^c} \quad \lambda > 0, c > 0$$

and cumulative distribution

$$(16) F(x) = 1 - e^{-\lambda x^c}$$

has a Weibull distribution, characterized by two parameters, λ and c . When $c = 1$ the distribution is exponential with parameter λ . [Mihram, pp. 279-281] attests to the utility of this distribution although rarely has it been proved to characterize actual processes.

Linear Regression and the Weibull Distribution

Suppose the n data points are assumed to be from a Weibull distribution with unknown parameters λ and c . The curve being sought is of the form

$$(17) y = F(x) = 1 - e^{-\lambda x^c}.$$

To employ a linear regression model, (17) is transposed and the logarithm twice taken to yield

$$(18) \ln(-\ln(1 - F(x))) = \ln(\lambda) + c \ln(x).$$

The right side of (18) is linear in $\ln(x)$. The linear regression model is applied to the pairs

$$(\ln(x_i), \ln(-\ln(1 - y_i))), \quad i = 1, \dots, n.$$

The m and b obtained satisfy

$$\ln(-\ln(1 - y)) = m \cdot \ln(x) + b$$

or

$$(19) \quad y = 1 - e^{-e^{\frac{b}{m}} x^m}.$$

Setting $c = m$ and $\lambda = e^{\frac{b}{m}}$ gives the distribution of (16).

Distributions of File Sizes at the Purdue University Computing Center

The disk file management system at Purdue is being modeled as part of a larger performance evaluation project. The initial part of the study tests the feasibility of a technique suggested by Forest Baskett. He collected the lengths of data files at the Stanford linear accelerator computing facility and attempted to fit them to a specific distribution.

At Purdue, several thousand files are maintained on a CDC 821 disk, and each file is one to (approximately) 600,000 (6 bit) bytes in length. A terminal-oriented system [Rosen et al.] and a batch system both access the disk (capacity: 32,678,000 bytes), which allocates space in multiples of sectors (1 sector = 640 bytes).

On April 22, 1973 the lengths of 2839 files (all files in the system) were measured. For each file, its length (in sectors), the number of days of its inactivity and the "type" of its user were recorded.*

See Table 1 for a summary of the length distribution.

* These items of data were selected because of their availability in an existing reporting procedure. User types include: student, faculty unsponsored research, faculty sponsored research, external user and system staff. Users do not purge their own files when they become inactive, and the system does it automatically after 15 days.

It was first assumed that file sizes were determined approximately by an exponential process. This was suggested by the memoryless property of the exponential: The probability of a file achieving length $L + B$ conditioned on the probability that the file has achieved length B is independent of the value B .

Attempts to discover an exponential process that matched the data were fruitless. When a distribution adequately matched short files, almost no lengthy files were predicted. When a distribution adequately matched lengthy files, very few short files were predicted.

As part of the attempt at exponential fitting, the natural logarithm of the complement of the cumulative distribution was plotted. Had the file sizes been exponentially distributed, the plot would have been linear. (See equation(8).) The nonlinear plot suggested a more general distribution. Using the methods above, a hyperexponential distribution was tested. Splitting the data into two groups (short and long files) and calculating w , λ_1 and λ_2 for both groups gave surprising results. The parameters obtained varied widely depending upon the choice of cutoff point. However the λ_1 and λ_2 differed by an order of magnitude (see Table 2) so that the fit for each group should have been approximately linear. Failing that, it appeared that the data were not generated as a sum of exponential distributions. Finally, the more general Weibull distribution was attempted. Using the method of the section above, parameters were obtained to yield the distribution

$$(20) \quad y = 1 - e^{-.34x^{.44}}.$$

The linear regression line was correlated .998 with the observed data.

In an attempt to validate equation (20), data were collected to replicate the experiment. On March 6, 1974 the lengths of all files were

again measured. Expanded use of the system resulted in 7944 files, an increase of almost 180% over the prior year volume. The distribution on March 6, 1974 was

$$(21) \quad y = 1 - e^{-.46x^{.41}},$$

with the same high correlation, .998.

Related Study and Acknowledgement

Forest Baskett proposed an incomplete gamma distribution to fit the data he collected. Attempts to duplicate his efforts failed. Herman Rubin suggested a model of file size generation assuming that each user has a fixed probability of terminating the file at each sector and that the user probabilities follow a beta distribution. John R. Rice gave insight into the least squares fitting technique. William Dahl and Richard Kovarik supplied the file length data.

References

- Freund, John E., Paul E. Livermore and Irwin Miller. Manual of Experimental Statistics, Prentice-Hall, Inc., Englewood Cliffs, N.J., 1960.
- Rosen, Saul, John M. Steele and R. E. Wagner. PROCSY -- The Purdue Remote On-line Console System. Proceedings of 1971 Annual Conference, Association for Computing Machinery, August 3-5, 1971, Chicago, Illinois, 253-262.
- Mihram, G. Arthur. Simulation, Statistical Foundations and Methodology, Academic Press, New York, 1972.

TABLE 1.
ACTUAL FREQUENCIES OF FILE LENGTHS (IN SECTORS)

LENGTH	FREQ.	LENGTH	FREQ.	LENGTH	FREQ.	LENGTH	FREQ.	LENGTH	FREQ.
1	755	50	2	101	2	162	1	259	2
2	332	51	7	102	4	163	5	267	2
3	157	52	8	103	4	164	1	271	1
4	102	53	5	104	1	165	1	273	1
5	123	54	4	105	4	167	2	274	1
6	85	55	8	106	7	168	3	276	1
7	84	56	7	107	5	169	1	278	1
8	62	57	3	108	2	170	1	283	1
9	59	58	5	109	1	171	1	285	1
10	44	59	10	110	2	174	1	287	1
11	30	60	5	111	3	175	2	295	1
12	34	61	4	112	3	176	2	303	1
13	37	62	5	113	2	178	2	305	1
14	28	63	2	114	5	179	2	309	1
15	33	64	2	115	2	180	2	312	1
16	27	65	6	116	2	181	1	324	1
17	18	66	3	118	3	183	1	341	1
18	23	67	8	119	2	185	2	342	1
19	17	68	4	120	1	188	3	346	2
20	22	70	2	121	1	190	1	362	1
21	20	71	2	122	2	192	3	369	3
22	21	72	5	123	2	194	1	372	1
23	17	73	3	124	1	198	1	388	2
24	6	74	2	125	2	200	1	395	1
25	13	75	3	126	1	201	1	396	1
26	15	76	4	129	2	202	1	404	1
27	10	77	6	130	1	203	2	407	1
28	10	78	2	131	2	205	1	414	1
29	18	79	2	133	2	206	1	423	1
30	17	80	2	136	3	207	1	424	1
31	10	81	1	137	1	208	1	453	1
32	18	82	3	138	1	210	1	458	1
33	18	83	2	141	1	212	2	472	1
34	14	84	2	142	1	216	1	517	1
35	10	85	2	143	1	217	1	545	1
36	11	86	6	144	1	219	1	551	1
37	10	87	3	145	1	220	1	564	1
38	9	88	1	147	2	223	1	566	1
39	10	89	3	149	2	224	1	576	1
40	17	91	2	150	1	225	2	580	1
41	21	92	3	151	1	229	1	581	1
42	3	93	3	153	1	230	1	594	1
43	6	94	6	154	1	232	2	596	1
44	6	95	2	155	6	237	1	612	1
45	6	96	1	156	1	239	1	695	1
46	9	97	15	157	1	248	2	716	1
47	7	98	4	158	1	252	1	910	1
48	7	99	3	160	2	253	1	1065	1
49	7	100	1	161	2	257	1		

TABLE 2.
EXPONENTIAL DISTRIBUTIONS APPROXIMATING OBSERVED DATA

<u>RANGE OF FILE LENGTHS OVER WHICH FIT WAS COMPUTED</u>	<u>PARAMETER OF EXPONENTIAL DISTRIBUTION</u>	<u>CORRELATION COEFFICIENT OF FIT</u>
1-5	.100	.98
1-10	.073	.98
1-15	.057	.97
1-20	.047	.97
1-25	.040	.97
1-30	.036	.96
1-50	.027	.97
5-1065	.008	.97
10-1065	.008	.98
15-1065	.008	.98
20-1065	.008	.98
25-1065	.008	.98
30-1065	.008	.98
50-1065	.008	.98
100-1065	.007	.99
200-1065	.006	.99
300-1065	.006	.98
1-1065	.008	.97